

1    **Are climate model simulations of clouds improving? An evaluation using the ISCCP**  
2    **simulator**

3    Stephen A. Klein<sup>1</sup>, Yuying Zhang<sup>1</sup>, Mark D. Zelinka<sup>1</sup>, Robert Pincus<sup>2</sup>, James Boyle<sup>1</sup>, and  
4    Peter J. Gleckler<sup>1</sup>

5    <sup>1</sup>Program for Climate Model Diagnosis and Intercomparison, Lawrence Livermore  
6    National Laboratory, Livermore, California, USA

7    <sup>2</sup>University of Colorado and NOAA/Earth System Research Laboratory, Boulder,  
8    Colorado, USA

9    Submitted to *Journal of Geophysical Research – Atmospheres*, July 2012; Revised, xxxx  
10   2012; Accepted, xxxx 2012

11   Corresponding author: S. A. Klein, Program for Climate Model Diagnosis and  
12   Intercomparison, Lawrence Livermore National Laboratory, 7000 East Avenue, L-103,  
13   Livermore, CA 94551. (klein21@llnl.gov)

14   **Running Title:** Evaluating Clouds in Climate Models

15   **Key Points**

- 16        • Newer climate models have improved simulations of cloud optical depth

17       • Cloud amount and cloud-top pressure simulations show smaller improvement

18       • Newer models have fewer compensating errors in their radiation budget

19   **Abstract**

20   The annual cycle climatology of cloud amount, cloud-top pressure and optical thickness  
21   in two climate model ensembles is compared to satellite observations in order to identify  
22   changes over time in the fidelity of climate model simulations of clouds. In more recent  
23   models, there is widespread reduction of a bias associated with too many highly reflective  
24   clouds, with the best models having eliminated this bias. With increased amounts of  
25   clouds with lesser reflectivity, recent models have reduced the compensating errors that  
26   permit models to simulate the time-mean radiation balance. Errors in cloud amount as a  
27   function of height or climate regime on average show little change or small improvement,  
28   although greater improvement can be found in the models of individual modeling centers.

29   **Index Terms:** 3337 Atmospheric Processes: Global climate models (1626, 4928); 3310  
30   Atmospheric Processes: Clouds and cloud feedbacks;       3360 Atmospheric Processes:  
31   Remote sensing (4337)

32   **Keywords:** clouds, climate models, satellite simulator

## 1. Measuring changes in the simulations of global cloudiness over time

The representation of clouds by climate models is a key ongoing challenge in the numerical representation of Earth's climate. Due to their large impact on Earth's radiation budget, clouds are important for determining aspects of current climate, such as surface air temperatures in many regions [*Ma et al.*, 1996; *Curry et al.*, 1996], the strength and variability of atmospheric circulations [*Slingo and Slingo*, 1988], and the magnitude of climate changes that result from perturbations in the chemical composition of the atmosphere [*IPCC*, 2007]. While important, the modeling of clouds is very difficult because most cloud processes happen at scales far smaller than can be resolved by climate models, and thus their bulk effects must be represented with imperfect parameterizations.

A large effort of many scientists over several decades and on several fronts has been undertaken to improve our understanding of cloud processes, often with the ultimate goal of improving the modeling of clouds in climate models. Observational programs have been launched to better understand cloud processes [*Stephens et al.*, 2002; *Ackerman and Stokes*, 2003], while very high-resolution models capable of resolving cloud processes provide additional information for the development of cloud parameterizations unavailable from observations [*GEWEX Cloud System Science Team*, 1993]. The community of scientists that work on physical process parameterizations in climate models has used the information provided by observations and fine-scale models to develop and implement many new and improved cloud parameterizations. Cloud

simulations in climate models may also be improved indirectly by complementary model development efforts that improve the representation of other physical processes including atmospheric dynamics, as well as by increases in model resolution.

Given this effort, it is important to ask: are climate model simulations of clouds improving and, if so, by how much? Here, we analyze the ability of two generations of climate models to simulate the climatological distribution of clouds and judge fidelity by comparison to several decade's worth of satellite observations. Because of the significant differences between the ways clouds are observed and the ways they are represented in models, we use a "satellite simulator" to increase the chances that differences between the models and observations represent actual model deficiencies. We find that significant progress in the ability of models to simulate clouds has occurred over the last decade, particularly in reducing the over-prediction of highly reflective clouds [Zhang *et al.*, 2005].

## **2. Climate Models, Satellite Observations, ISCCP Simulator and Analysis Methods**

### **2.1 Climate Models**

The models we analyze are those that submitted output to the first two phases of the Cloud Feedback Model Intercomparison Project [McAvaney and LeTreut, 2003; Bony *et al.*, 2011]. Submissions to the first phase (CFMIP1) were completed by the end of 2005 and thus the nine models (Table 1) we analyze were all formulated prior to that

time, with HadSM3 being perhaps the oldest of these models. Submissions to the second phase (CFMIP2) began in late 2011 and as of the time of this writing<sup>1</sup> we have output from eight models (Table 2). CFMIP2 is a subset of the much wider fifth Coupled Model Intercomparison Project (CMIP5) [Taylor *et al.*, 2012] associated with the fifth assessment report of the Intergovernmental Panel on Climate Change. Although less formal, there was also a close connection between CFMIP1 and the corresponding third Coupled Model Intercomparison Project (CMIP3) [Meehl *et al.*, 2007]. As some models that participated in CFMIP1 did not participate in CMIP3, we retain the more accurate label of CFMIP (instead of CMIP) when referring to the ensembles.

A direct evaluation of model changes is complicated by the fact that the CFMIP1 output used here is from the control climate integrations of slab-ocean models (i.e., atmospheric models coupled with a mixed-layer model of the upper ocean), while the CFMIP2 output is from simulations of the atmosphere model with sea surface temperatures and sea-ice distributions prescribed from observations from recent decades (i.e. Atmospheric Model Intercomparison Project (AMIP) simulations [Gates *et al.*, 1999]). This difference arises because the satellite simulator output we require is only available from the slab-ocean models of CFMIP1, while the slab-ocean model framework is not part of CFMIP2. Nonetheless, we believe that the difference in modeling framework has only a minor

---

<sup>1</sup> We intend to add other CFMIP2 models to our analysis should they become available during the review process. We think it is possible to add results from the GFDL AM3 and EC-Earth models.

impact, because the differences in surface boundary conditions between slab-ocean models and AMIP integrations (and hence the resulting distribution of clouds) are small, even for slab-ocean models constructed to mimic the climate of the pre-industrial era. We have tested this notion by comparing AMIP and slab-ocean model simulations for one model (CAM4), and find that differences in our results resulting from the different modeling frameworks to be much smaller than differences among CFMIP models.

## 2.2 Satellite Observations

We compare the clouds simulated by climate models to the cloud climatology of observations created by the International Satellite Cloud Climatology Project (ISCCP) [Rossow and Schiffer, 1991, 1999]. ISCCP provides estimates of the area coverage of clouds stratified by *ctp*, the apparent cloud-top pressure of the highest cloud in a column, and by  $\tau$ , the column integrated optical thickness of clouds. These estimates are the results of retrieval algorithms applied to radiance observations from the visible and infrared window channels of geostationary and polar orbiting satellites. They are accumulated for 280 km x 280 km regions every 3 hours starting in July 1983 and we use data through June 2008. Area coverage estimates are summarized in a joint histogram with 6 bins in  $\tau$  and 7 bins in *ctp*; bin boundaries are shown in Figure 7. We use custom-built daytime-only monthly averages that are described more fully in Pincus et al. [2012] and are available from <http://climserv.ipsl.polytechnique.fr/>.

As a point of comparison, we also use roughly analogous observations from the

MODerate Resolution Imaging Spectrometer (MODIS) instruments for the period March 2000 through April 2011 [*Pincus et al.*, 2012]. MODIS uses substantially different methods of estimating *ctp* than does ISCCP, so the amounts of clouds in each bin of the joint histogram of *ctp* and  $\tau$  from MODIS are not comparable to those observed by ISCCP or the output of an ISCCP simulator applied to climate models. (MODIS observations may be compared to the output of a MODIS simulator [*Pincus et al.*, 2012], but that was not available at the time of CFMIP1.) On the other hand, MODIS retrievals of  $\tau$  are roughly equivalent to those from ISCCP, so we compare MODIS observations, aggregated over bins of *ctp*, to both ISCCP observations and the output of ISCCP simulators.

## 2.3 ISCCP Simulator

A satellite simulator is a diagnostic code applied to model variables that reduces the influences of inconsistencies between the ways clouds are observed and the ways they are modeled [*Bodas-Salcedo et al.*, 2011]. By mimicking the observational process in a simplified way, the simulator attempts to compute what a satellite would retrieve if the real-world atmosphere had the clouds of the model. Simulators increase the chances that the comparison of satellite retrievals to model output after run through a simulator is an evaluation of the fidelity of a model's simulation rather than a reflection of observational limitations or artifacts. The use of a satellite simulator also puts model intercomparison on a firmer basis by minimizing the impacts of how clouds are defined in different

131 parameterizations.

132 The ISCCP simulator is the oldest of the satellite simulators used to evaluate clouds in  
133 models and has been widely used by most major climate modeling centers since its  
134 creation over ten years ago [*Klein and Jakob, 1999; Webb et al., 2001*]. The ISCCP  
135 simulator mimics the key aspect of the ISCCP retrieval algorithms that radiances in every  
136 cloudy satellite pixel are assumed to arise from a single homogenous layer of cloud with  
137 *ctp* determined from an infrared brightness temperature. In detail, the ISCCP simulator  
138 takes a model's vertical profile of grid-box mean clouds and creates a set of sub-grid  
139 scale columns which are completely clear or cloudy at each level and which are  
140 consistent with the model's cloud-overlap parameterization. (This step is bypassed for  
141 models that provide to the simulator a set of previously generated sub-grid scale  
142 columns.) From every sub-grid scale column, one determines the single value of *ctp* and  
143 column-integrated  $\tau$  that would be consistent with the single-layer cloud retrieval that  
144 ISCCP applies to every cloudy satellite pixel. In this step, *ctp* is determined by applying a  
145 simplified radiative transfer model in each sub-grid scale column to determine an infrared  
146 brightness temperature, which is then converted to the temperature at cloud-top by using  
147 a cloud longwave emissivity derived from  $\tau$ , as in the ISCCP retrieval algorithm. Once a  
148 cloud-top temperature has been determined, *ctp* is equated with the interpolated pressure  
149 that has the identical temperature according to the model's profile of temperature. The  
150 column-integrated value of  $\tau$  is equated with the sum of model-reported  $\tau$  from all model  
151 layers that are cloudy in a given sub-grid scale column. From these sub-grid scale values



152 of  $ctp$  and  $\tau$ , the grid-box mean joint histogram of  $ctp$  and  $\tau$  is formed for every grid  
153 box and then subsequently averaged over time. To make the comparison with satellite  
154 retrievals of  $\tau$  more fair, the ISCCP simulator is only applied to grid-boxes that are sunlit  
155 at a given model time.

156 The ISCCP simulator itself changed between CFMIP1, which used v3.5, and CFMIP2,  
157 which used v4.1, raising the possibility that differences in the diagnostics might be  
158 mistaken for changes in simulation quality. The most significant algorithmic difference  
159 between these two versions involves the determination of  $ctp$  for clouds under  
160 atmospheric temperature inversions, such as subtropical marine stratocumulus. In these  
161 situations, ISCCP often erroneously assigns  $ctp$  to a level far higher (100 – 300 hPa) in  
162 the atmosphere than it should be [Garay *et al.*, 2008]. In CFMIP1,  $ctp$  is assigned to the  
163 highest interpolated pressure (lowest altitude) with matching cloud-top temperature, but,  
164 since the simulator is intended to mimic the retrieval process (even when it is faulty), the  
165 simulator was changed so that  $ctp$  is assigned to the lowest interpolated pressure (highest  
166 altitude) with matching cloud-top temperature when a temperature inversion is present in  
167 the model. We have verified that the impact of this and other simulator differences have  
168 little impact on our results by comparing the output of these two versions of the ISCCP  
169 simulator when applied to identical integrations of two CFMIP2 models (CAM4 and  
170 HadGEM2) (not shown). Simulator changes primarily affect  $ctp$  with differences of up to  
171 0.01 in the amounts of clouds annually averaged over the domain 60°N-60°S for  $ctp$  bins  
172 where  $ctp < 680$  hPa, and somewhat larger differences of up to 0.04 for  $ctp$  bins where

173  $ctp > 680$  hPa.

174 We only use models for which we are reasonably confident of a correct implementation  
175 of the ISCCP simulator. Our primary test is to verify that the sum of cloud cover over all  
176 bins of the joint histogram is consistent with the model diagnostic of total cloud cover  
177 ('clt') which a model computes without using the ISCCP simulator [Zelinka *et al.*, 2012].

## 178 2.4 Analysis Methods

179 Climatological joint histograms of  $ctp$  and  $\tau$  are formed for every calendar month by  
180 averaging model and observational data on a common  $2^\circ$  latitude by  $2.5^\circ$  longitude grid  
181 from every available year. Most model climatologies are based upon either 20 or 30  
182 simulated years whereas the observed climatologies are for 25 years for ISCCP and 11  
183 years for MODIS, but differences in the number of years available do not materially  
184 affect our evaluation [Pincus *et al.*, 2008]. (The scalar measures of the fidelity of model  
185 simulations [Section 4] are sensitive to this issue if the number of years used to form a  
186 climatology is very low ( $< 5$ ); this only affects results for the two MIROC models in  
187 CFMIP1.) To minimize issues with cloud retrievals above surfaces with snow or ice, we  
188 restrict our analysis to the domain  $60^\circ\text{N}$ - $60^\circ\text{S}$ .

189 We evaluate changes over time in two ways. One considers changes in the multi-model  
190 mean from each of the CFMIP ensembles. This has the advantage of considering all  
191 available models and of highlighting common model errors. However, multi-

model means are sensitive to the addition of new models (especially given the small sizes of the model ensembles) and changes in the multi-model mean may not reveal individual model error reductions when the spread of model results is centered on the observed value, as is often the case [Gleckler *et al.*, 2008]. To address these limitations, we also track the changes over time in the models from the four modeling centers that have contributed one or more models to both ensembles. For this analysis, we use models from the Canadian Centre for Climate Modeling and Analysis (AGCM4.0 to CanAM4), the United Kingdom’s Met Office Hadley Centre (HadSM3 to HadSM4 to HadGEM1 to HadGEM2), the Japanese climate model effort associated with MIROC (MIROC(hisens) and MIROC(losens) to MIROC5), and the United States climate modeling effort associated with the Community Atmosphere Model (CCSM3.0 to CAM4 to CAM5).

### **3. Comparisons of climate model simulations of clouds to satellite observations**

#### **3.1 Common improvements and failures in the simulation of total cloud amount**

We begin our analysis by examining the ability of models to simulate the space-time distribution of total cloud amount, i.e., how often a cloud occurs with any value of  $ctp$  and  $\tau$ , which is perhaps the most fundamental aspect of a model’s ability to simulate clouds. Unfortunately, this quantity is problematic to define from observations: satellite estimates of total cloud amount are extremely sensitive to many observational factors including the scale and sensitivity of the fundamental observations, as well as decisions made during the aggregation to larger scales [Stubenrauch *et al.*, 2009; Mace *et*

212 *al.*, 2009; *Marchand et al.*, 2010; *Pincus et al.*, 2012]. We make the comparison more  
213 robust by restricting the analysis to clouds with  $\tau$  exceeding some minimum threshold  
214  $\tau_{\min}$ , which we set to minimize hard-to-detect and partly-cloudy observations. We select  
215  $\tau_{\min} = 1.3$  from among the discrete choices offered by the bin boundaries of the joint  
216 histogram of  $ctp$  and  $\tau$  by balancing the following desires: (a) to maximize the number of  
217 clouds that we examine, (b) so that the observational datasets we use agree among  
218 themselves, ensuring robust model evaluation, and (c) to minimize the chances that an  
219 observational platform would have missed a cloud with  $\tau > \tau_{\min}$ . Setting  $\tau_{\min} = 1.3$   
220 provides the smallest relative bias and relative root-mean-square difference, as well as the  
221 maximum correlation coefficient, between the space-time distributions of the annual  
222 cycle climatologies of ISCCP and MODIS.

223 Figure 1 illustrates the annual mean total cloud amount for the multi-model means of the  
224 CFMIP1 and CFMIP2 ensembles, the ISCCP and MODIS observations, and the  
225 difference of the CFMIP2 multi-model mean with ISCCP observations and with the  
226 CFMIP1 multi-model mean. For the domain 60°N-60°S, the annual mean total cloud  
227 amount fraction with a  $\tau_{\min}$  of 1.3 from ISCCP and MODIS is 0.51 and 0.47, respectively.  
228 The multi-model means of both CFMIP1 and CFMIP2 are 0.43 with more than  $\frac{3}{4}$  of  
229 models in both ensembles below the range of observational estimates. Although the  
230 multi-model mean is identical between the two ensembles, if one examines these area-  
231 averaged values for the four model families in which we can track progress, in every case  
232 the most recent model is closer to the observational estimates. The increase is quite

233 striking for the Hadley Centre models, with HadSM3 having a total cloud amount of  
234 0.33 but HadGEM2 having a total cloud amount of 0.43.

235 Relative to ISCCP observations, model underestimates of total cloud amount  
236 preferentially occur in regions of subtropical marine stratocumulus on the eastern sides of  
237 subtropical ocean basins and over middle latitudes. In the stratocumulus regions, there is  
238 a wide variety of results in both ensembles with about 3-4 members in each ensemble  
239 having total cloud amount values close to observed and the remainder of models  
240 significantly below observational estimates. Although the differences between the multi-  
241 model means of ensembles are small in these regions, one finds marked progress in 3 out  
242 of the 4 families we can track with the amount of clouds in the most recent model  
243 versions close to observed. This suggests that at least for the modeling centers for which  
244 we can track progress, the simulation of current climate amounts of subtropical  
245 stratocumulus has been improving, perhaps in response to the well-known importance of  
246 the low clouds in these regions for mean climate and climate sensitivity [*Bony and*  
247 *duFresne*, 2005].

248 Although not as well known, models also typically underestimate total cloud amount at  
249 middle latitudes over both land and ocean (Figure 1). While a few models are close to  
250 observed over the middle latitude oceans, all models underestimate total cloud amount  
251 over the middle latitudes of Eurasia and North America. Examination of level-by-level  
252 cloud amount indicates that these underestimates, over both land and ocean, are primarily  
253 of lower level clouds ( $ctp > 560$  hPa) although underestimates in upper level

clouds ( $ctp < 560$  hPa) do contribute some to this error, depending on the model. When examining results by model families, one finds no consistent sign of progress for this bias either over ocean or land, consistent with larger middle-latitude bias in the CFMIP2 multi-model mean relative to the CFMIP1 multi-model mean.

### 3.2 Improvements as a function of cloud-top pressure and cloud optical depth

In addition to getting clouds to occur in the right places and times, having a good simulation of  $ctp$  and  $\tau$  is essential to getting the correct long- and shortwave impacts of a given cloud on the top-of-atmosphere radiation budget. Figure 2 illustrates the amount of clouds with  $\tau > 1.3$  as a function of  $ctp$  averaged over  $60^{\circ}\text{N}$ - $60^{\circ}\text{S}$ . Models tend to underestimate the amount of middle ( $440 \text{ hPa} < ctp < 680 \text{ hPa}$ ) and low-level ( $ctp > 680 \text{ hPa}$ ) clouds while having about the right amount of high-level ( $ctp < 440 \text{ hPa}$ ) clouds [Zhang *et al.*, 2005]. The general underestimate of low-level clouds is consistent with the lack of clouds in marine stratocumulus and middle-latitudes mentioned above. Differences in middle-level clouds are somewhat hard to interpret as many middle-level clouds observed by ISCCP are in fact multi-layer cloud scenes of cirrus above boundary layer cloud [Marchand *et al.*, 2010; Mace *et al.*, 2011]. Though the ISCCP simulator is capable of reproducing this artifact [Mace *et al.*, 2011], it will do so only if a model produces thin cirrus over boundary layer clouds. Thus, underestimates of middle-level cloud may actually indicate a lack of cirrus above boundary layer cloud.

Relative to that of the CFMIP1 ensemble, the CFMIP2 multi-model mean is closer to

the observed amounts for 6 out of 7 bins of *ctp*, suggesting some improvement. This improvement is noticeable in the relative amounts of low-level clouds in the two lowest *ctp* bins. While a large part of this improvement is due to the change in the simulator's determination of *ctp* for clouds under an inversion, improvement can be found in the models from modeling centers that contribute more than one model to a given ensemble (compare HadSM3 to HadGSM1 and CAM4 to CAM5). Because the ISCCP simulator version does not change within these two pairs, we can conclude that these models have improved their simulation of low-level clouds. For middle-level clouds, there is also a reduction in the model underestimate, particularly for the 560-680 hPa *ctp* bin. In fact, the perfect agreement of CAM5 with ISCCP for this bin can be attributed to the fact that snow is now radiatively active and thus the simulator counts the contribution of snow to  $\tau$  and the infrared-brightness temperature used to determine *ctp* [Kay *et al.*, 2012].

Figure 3 illustrates the amount of clouds as a function of  $\tau$  regardless of *ctp* and averaged over 60°N-60°S. More so than in the case of *ctp*, rather marked improvement can be seen for  $\tau$  bins where ISCCP and MODIS agree fairly well ( $\tau > 3.6$ ). In particular, the amounts of optically thick clouds ( $\tau > 23$ ) are significantly closer to observed in the CFMIP2 ensemble relative to the CFMIP1 ensemble with a marked reduction in the previously identified overestimate of highly reflective clouds [Zhang *et al.*, 2005]. All but one of the CFMIP2 models have fewer clouds in the optically thickest bin ( $\tau > 60$ ) than all but one of the CFMIP1 models. This bias reduction is widespread enough that it is dramatically

present for each of the 4 model families in which we can track progress (Figure 4).

The fraction of the 60°N-60°S area covered by optically thick cloud is 0.175 for the CFMIP1 ensemble mean but is 0.121 for the CFMIP2 ensemble mean. The CFMIP2 ensemble mean is still larger than the observational estimates of 0.064 for ISCCP and 0.082 for MODIS, indicating that about half of the bias remains. For HadGEM2 and MRI-CGCM3, the amount of optically thick cloud is within the range of the two observational estimates. The reduction between ensembles in optically thick clouds is larger for lower-level ( $ctp > 560$  hPa) clouds than it is for upper-level ( $ctp < 560$  hPa) clouds, 0.043 vs. 0.009 respectively. With the greater reduction in lower-level optically thick clouds, 7 out of 8 CFMIP2 models as opposed to 5 out of 9 CFMIP1 models reproduce the fact that optically thick clouds occur more frequently with  $ctp$  at upper levels than at lower levels. However, for only 2 CFMIP1 and 3 CFMIP2 models does the ratio of upper to lower-level optically thick clouds exceed the observed value of 1.7 for ISCCP (2.2 for MODIS).

Geographically, one can see from the multi-model means that the significant reductions in the amount of optically thick clouds occur over both the subtropical stratocumulus regions and middle-latitude land and especially ocean (Figure 5). There is no improvement in the multi-model mean overestimate of optically thick clouds over tropical continents, and this bias is present in 7 out of 9 CFMIP1 models and 7 out of 8 CFMIP2 models. We suspect that the common model bias in the diurnal cycle precipitation over tropical land [Yang and Slingo, 2001; Dai, 2006] contributes to this



315 error by producing too many optically thick anvil clouds near mid-day, when they are  
316 visible to the ISCCP simulator, rather than at night.

317 The decrease in optically thick clouds has been accompanied by an increase in the  
318 amount of clouds with intermediate optical depths ( $3.6 < \tau < 23$ ) (Figures 3 and 6). This  
319 increase is present in each of the 4 model families for which we can track progress, with  
320 one 1 CFMIP model (IPSL-CM4) and 3 CFMIP2 models (CAM5, HadGEM2, MPI-  
321 ESM-LR) having the amount of intermediate optical depth clouds lying in between the  
322 values from ISCCP and MODIS.

323 Passive observational estimates of the amount of cloud with  $0.3 < \tau < 3.6$  disagree  
324 sharply, in part because many of the observations which produce clouds in this optical  
325 thickness range are partly cloudy [*Pincus et al.*, 2012]. This makes it impossible to assess  
326 the fidelity of model simulations for these clouds. For  $\tau < 0.3$ , there is a wide variety of  
327 model results, particularly in CFMIP1 where the two MIROC models each have more  
328 than 0.25 of the area covered by clouds of this optical depth range. Clouds this thin have  
329 too little contrast on the top-of-atmosphere radiation budget to be detected with the  
330 passive sensors used by ISCCP and MODIS; in fact, the  $\tau$  bin boundary of 0.3 is chosen  
331 to crudely mimic a sensitivity threshold for ISCCP (W. B. Rossow, personal  
332 communication). Assessment of very thin clouds requires the use of an active sensor such  
333 as CALIPSO [*Winker et al.*, 2009]. Such an assessment would be relevant for the  
334 plentiful but very thin tropopause-level cirrus in the tropics [*Mace et al.*, 2009; *Thorsen et*

335 *al.*, 2011].

### 336 3.3 Radiative impact of model errors in cloud properties

337 As in nature, clouds in climate models strongly affect the radiation balance as a function  
338 of space and time. Model tuning guarantees that the global and annual average of the net  
339 radiation is close to zero, but significant regional errors in the radiation field may persist,  
340 and correct regional fluxes can be achieved through compensating errors in cloud  
341 properties. One common error is to have clouds which are too few but too bright, that is,  
342 to have lower-than-observed cloud amounts with larger-than-observed values of  $\tau$ , such  
343 that the average shortwave radiation budget is about right [Zhang *et al.*, 2005].

344 We explore these issues by using cloud radiative kernels [Zelinka *et al.*, 2012] to compute  
345 the radiative effects of errors in cloud properties. A cloud kernel  $K^{SW,LW}$  is the result of a  
346 radiative transfer calculation that computes the impact on the top-of-atmosphere short-  
347 and long-wave fluxes, relative to clear-sky, of the addition of a unit area covered by a  
348 cloud with a given *ctp* and  $\tau$ . Our kernels are computed as a function of latitude,  
349 longitude and calendar month. Multiplying the kernels by the bias, relative to ISCCP, in  
350 cloud amount in each bin of the joint histogram yields an estimate of the error in top-of-  
351 atmosphere radiation budget due to errors in the simulated distribution of clouds as a  
352 function of *ctp* and  $\tau$ .

353 Figure 7 shows the annually and 60°N- 60°S averaged bias relative to ISCCP in

cloud amount fraction in the joint histograms of  $ctp$  and  $\tau$  for the four model families in which we can track progress and the multi-model means for CFMIP1 and CFMIP2. Figure 8 and 9 show the corresponding biases in  $\text{W m}^{-2}$  for the short- and long-wave radiation of the same models. (The Canadian model pairing is absent from Figures 8-9 because we cannot perform accurate cloud kernel calculations for AGCM4.0 for the reasons discussed in the Appendix of *Zelinka et al.* [2012].) The oldest models are in the left column and the most recent models on the right. The prominent overestimate of optically thick clouds occurs in nearly all  $ctp$  bins in the earlier models, but is much reduced in the more recent set. Likewise the underestimate of optically thin ( $0.3 < \tau < 3.6$ ) and intermediate clouds present in nearly all  $ctp$  bins has been reduced in the more recent model versions. As discussed above, whether or not the biases in thin clouds are real is unclear.

The radiative impact of these biases on the short-wave spectrum quantifies the nature of compensating errors (Figure 8), with the overestimates of reflected shortwave by clouds with  $\tau > 23$  compensating for a lack of reflection by clouds with thin and intermediate optical depths. The figure is similar to that of the cloud biases (Figure 7) except that weighting by the shortwave radiative kernel reduces the impact of the underestimate of optically thin clouds relative to the overestimate of optically thick clouds. The degree of compensation is markedly reduced in the more recent models. For example, in HadSM3 there was  $27 \text{ W m}^{-2}$  too much reflectance by clouds with  $\tau > 60$ , whereas in the most recent model HadGEM2-A, the bias is less than  $1 \text{ W m}^{-2}$ . Similarly, for HadSM3,

CCSM3.0 and MIROC (hisens) and MIROC (losens), there was an underestimate of reflected shortwave radiation by clouds with  $3.6 < \tau < 9.4$  of about  $10 \text{ W m}^{-2}$ , but in CAM5, MIROC5, and HadGEM2-A this bias is less than  $3 \text{ W m}^{-2}$ . In the multi-model mean, too much reflectance by optically thick clouds compensates for an underestimate in reflection by clouds with  $1.3 < \tau < 9.4$  for most *ctp* bins, but the biases are smaller in the more recent models.

In the longwave spectrum, the nature of compensating biases is similar but with emphasis on upper level clouds (Figure 9). In general, there is too much reduction of outgoing longwave radiation by high clouds with  $\tau > 60$ , which compensates for a lack of reduction of outgoing longwave radiation by thinner clouds at both middle and high levels of the troposphere. The progress is clearly identifiable but not quite as prominent as in the case of shortwave radiation with noticeable progress for the Community Atmosphere and Hadley Centre models but less so for the MIROC model and the multi-model means.

#### 4. Scalar measures of the fidelity of model simulations

While the evidence just presented supports the notion that the simulation of clouds in climate models has been improving, it is helpful to provide scalar measures of the fidelity of model simulations that can quantitatively demonstrate progress. Here we present a few such quantities chosen to measure different aspects of cloud simulations and for which

394 observational uncertainty is less than the differences between models and observations  
 395 and among models themselves. These measures might be considered for a list of metrics  
 396 for clouds in climate models [Gleckler *et al.*, 2008; Pincus *et al.*, 2008; Williams and  
 397 Webb, 2009], although we do not develop this aspect here.

398 In the following,  $c(ctp, \tau, X)$  is the amount of cloud in a given bin of the ISCCP  
 399 histogram and is a function of cloud-top pressure  $ctp$ , optical depth  $\tau$ , latitude, and  
 400 generalized position  $X$ , including latitude, longitude, and month. Total cloud amount  
 401  $C(\tau_{\min})$  is the sum of the cloud amounts of all bins with  $\tau$  greater than the minimum  
 402 optical thickness  $\tau_{\min}$ :

$$403 \quad C(\tau_{\min}, X) = \sum_{ctp} \sum_{\tau}^{\tau > \tau_{\min}} c(ctp, \tau, X) \quad (1)$$

404 We compute the normalized root-mean-square error  $Z_1$  in the space-time distribution of  
 405 total cloud amount, as:

$$406 \quad Z_1(\tau_{\min}) = \sqrt{\int_X [C^{MOD}(\tau_{\min}, X) - C^{OBS}(\tau_{\min}, X)]^2} / \sigma_1. \quad (2)$$

407 The integral in (2) denotes the area-weighted space-time average of squared differences  
 408 between the model and ISCCP observations. The root-mean-square differences are  
 409 normalized by the space-time standard deviation of the observed total cloud amount,

410 given by:

$$411 \quad \sigma_1 = \sqrt{\int_X \left[ C^{OBS}(\tau_{\min}, X) - \bar{C}^{OBS}(\tau_{\min}) \right]^2}. \quad (3)$$

412 As in Section 3.1, we set  $\tau_{\min} = 1.3$ .

413 Equation (1) uses the ISCCP simulator to ensure that model definitions of cloudiness are  
 414 comparable with what is robustly observable but ignores the wealth of information  
 415 provided by the joint histogram of  $ctp$  and  $\tau$ . We evaluate the error  $Z_2$  in this more finely-  
 416 resolved distribution as the sum over a finite number of cloud-top pressure ( $N_{ctp}$ ) and  
 417 optical thickness ( $N_\tau$ ) bins of squared differences between the model and ISCCP  
 418 observations:

$$419 \quad Z_2 = \sqrt{\int_X \frac{1}{N_{ctp} \times N_\tau} \times \sum_{ctp} \sum_{\tau}^{\tau > \tau_{\min}} \left( c^{MOD}(ctp, \tau, X) - c^{OBS}(ctp, \tau, X) \right)^2} / \sigma_2. \quad (4)$$

420 This measure is sensitive to differences in each bin with  $\tau > \tau_{\min}$ , and would be  
 421 applicable if the ISCCP simulator were capable of reproducing every aspect of the ISCCP  
 422 observational processes. But comparisons with clouds retrieved from ground-based  
 423 remote sensors and passed through the ISCCP simulator [Figures 2c and 3c of *Mace et*  
 424 *al.*, 2011] suggest that the accuracy of ISCCP retrievals is about  $\pm 200$  hPa for  $ctp$  and a  
 425 factor of 3 for  $\tau$ . We therefore compute  $Z_2$  from a reduced-resolution histogram with bin  
 426 boundaries in  $ctp$  of 440 hPa and 680 hPa and in  $\tau$  of 3.6 and 23. (This is equivalent

to the reduced-resolution joint histogram available in the monthly-averaged ISCCP data archives.) Considering the greater uncertainty of thin-cloud retrievals, we set  $\tau_{\min} = 3.6$ , and calculate differences only for the 6 bins with  $\tau > \tau_{\min}$ .  $Z_2$  is normalized by  $\sigma_2$ , the accumulated space-time standard deviation of observed cloud amounts in the reduced bin set, making  $Z_2$  the normalized root-mean-square error in the amount of low-intermediate, low-thick, medium-intermediate, medium-thick, high-intermediate, and high-thick clouds.

We compute the radiatively-relevant error  $Z_3$  in the distribution of clouds by using the radiative kernels to weight bin-by-bin errors by their radiative impact on top-of-atmosphere radiation fluxes:

$$Z_3^{SW,LW}(\tau_{\min}) = \sqrt{\int_X \frac{1}{N_{ctp} \times N_{\tau}} \times \sum_{ctp} \sum_{\tau}^{\tau > \tau_{\min}} \left[ K^{SW,LW}(ctp, \tau, X) \times (c^{MOD}(ctp, \tau, X) - c^{OBS}(ctp, \tau, X)) \right]^2} / \sigma_3^{SW,LW} \quad (5)$$

Multiplication by radiative kernel is performed for each bin of the original ISCCP histogram before aggregation to the reduced bin set. This measure  $Z_3$  has separate components for the shortwave and longwave spectrum, and is normalized by the accumulated space-time standard deviation of the radiative impacts of observed clouds from the reduced bin set.

Figure 10 shows  $Z_1$ ,  $Z_2$ ,  $Z_3^{LW}$ , and  $Z_3^{SW}$  for each model stratified into two rows according to the model ensemble. Arrows from earlier to later models indicate the change with time in the fidelity of model simulations; left-pointing arrows indicate smaller errors over time. The arrows connect the earliest and latest models from the modeling centers for which we track progress as well as the mean measure of each model ensemble. (In order to identify progress over time, the mean only includes the earliest CFMIP1 (latest CFMIP2) models from modeling centers that contribute more than one model to a given ensemble.)

For the total cloud amount measure  $Z_1$ , values range from 0.65 to 1.18 indicating that the standard deviation of biases in total cloud amount relative to ISCCP are generally comparable in size to the space-time of standard deviation of observed total cloud amount. To put this number into context, the  $Z_1$  measure between the MODIS and ISCCP climatologies is 0.47. All model differences with ISCCP exceed this value, so it is likely that errors in the climatology of total cloud amount are robustly determined. Consistent with Figure 1, there is not a clear sign of improvement when considering the ensemble as a whole with the CFMIP1 ensemble mean value of  $Z_1$  equal to 0.86 and the CFMIP2 ensemble mean value of  $Z_1$  equal to 0.82. However, improvement is found for the Hadley Centre and Community Atmosphere models with a reduction of  $Z_1$  from 1.12 for HadSM3 to 0.70 for HadGEM2A and a reduction of  $Z_1$  from 0.94 for CCSM3.0 to 0.65 for CAM5, with little change in  $Z_1$  for the Canadian and MIROC models or the ensemble mean.



466 For the cloud property measure  $Z_2$ , much more dramatic progress can be found. For  
 467 three of the 4 models in which we can track progress (Hadley Centre, Community  
 468 Atmosphere, and Canadian Centre models), errors relative to ISCCP has been reduced by  
 469 40-45% (relative), from 150-175% to 80-105% of the standard deviation of the ISCCP  
 470 amounts of the 6 intermediate and thick cloud types. For the ensemble mean measure,  
 471 more moderate progress can be found with 15-30% (relative) reduction in  $Z_2$ . Separate  
 472 calculations reveal that the majority of the improvement in  $Z_2$  comes from a better  
 473 simulation of the amounts of optically intermediate ( $3.6 < \tau < 23$ ) and thick ( $\tau > 23$ )  
 474 clouds, than it does for improvements in the high, middle, and low amounts of clouds  
 475 (with  $\tau > 3.6$ ) (figures not shown). For the equivalent error measure calculated using only  
 476 two bins for optically intermediate and thick clouds regardless of  $ctp$ , the value for the  
 477 best model HadGEM2A is close to that calculated for differences between the observed  
 478 ISCCP and MODIS distributions (0.70 vs. 0.59).

479 Radiatively-relevant cloud property measures  $Z_3^{SW}$  and  $Z_3^{LW}$  are shown in the bottom row  
 480 of Figure 10. Similar to the cloud property measure  $Z_2$ , both measures show significant  
 481 error reductions of 20-30% for the ensemble mean measure with larger 40-50% error  
 482 reductions for individual models such as those of the Hadley Centre and Community  
 483 Atmosphere. Again, the majority of this error reduction comes from improvement in the  
 484 simulation of  $\tau$ , indicating that models are better simulating the amount of shortwave  
 485 radiation reflected and longwave radiation trapped by optically intermediate and thick  
 486 clouds. Though it may appear that there is a redundancy among  $Z_2$ ,  $Z_3^{SW}$  and  $Z_3^{LW}$ , only

487  $Z_2$  and  $Z_3^{\text{SW}}$  are highly correlated; all other possible pairings, including those with  $Z_1$ ,  
488 have statistically insignificant inter-model correlations.

489 **5. Why are simulations of clouds improving, and what impacts might this have?**

490 The agreement between satellite observations and simulations by climate models of the  
491 climatological annual cycle of cloud amount, cloud-top pressure, and optical thickness  
492 has improved over the last decade. The improvement is most striking in the simulation of  
493  $\tau$ , where a bias of having too many optically thick clouds ( $\tau > 23$ ) has been reduced by  
494 about 50% in the multi-model mean, with the best models having eliminated this bias.  
495 With a corresponding increase in the simulated amount of clouds with intermediate  
496 optical depth ( $3.6 < \tau < 23$ ), this reduces the tendency for climate models to simulate  
497 approximately the right amount of shortwave radiation reflected by clouds but with the  
498 compensating errors of having too few clouds that are too bright.

499 Improvement in the amount or height distribution of clouds is not clear in the ensemble  
500 as a whole although progress can be found in individual models. For example, the  
501 simulations of total cloud amount in the Hadley Centre and Community Atmosphere  
502 models do show noticeable improvement (see  $Z_1$  of Figure 10); in part, this improvement  
503 results from better simulations of the amount of clouds in the climatically important  
504 subtropical marine stratocumulus regions, where the amount of cloud is close to that  
505 observed in their most recent models. Some things show no improvement in the majority  
506 of climate models such as the underestimate of cloud over middle-latitudes, particularly

507 over land, and an overestimate in the amount of optically thick cloud over tropical land.

508 Pinpointing the reasons for model improvement is difficult without testing of individual  
509 modifications from among the myriad of changes that modeling centers have  
510 implemented in the last decade, and it is likely that many factors have contributed. Even  
511 apart from parameterization changes, the incorporation of ISCCP simulator diagnostics in  
512 the routine evaluation of developmental model versions (as was done at the Hadley  
513 Centre for much of the last decade [*Martin et al.*, 2006]) can have a subtle but persistent  
514 influence on the choices made in the model-development process in such a way as to  
515 lead to improved simulation of clouds.

516 With regard to parameterizations, the improved boundary layer turbulence and shallow  
517 convection parameterizations in the Hadley Centre and Community Atmosphere models  
518 [*Lock et al.*, 2000; *Bretherton and Park*, 2009; *Park and Bretherton*, 2009] are almost  
519 certainly responsible for the improved simulations in marine stratocumulus clouds.  
520 However, in the case of the improved optical depth distribution, the causes of  
521 improvement are less clear but there are some clues from what has happened at the  
522 individual modeling centers whose progress we can track.

523 Beginning with the Canadian model, the reduction in the amount of optically thick cloud  
524 between its two versions is striking given the relatively few changes between model  
525 versions (J. Cole, personal communication). The likeliest cause is thought to be the  
526 introduction into CanAM4 of a new treatment of sub-grid scale variability in

cloud optical properties known as the Monte Carlo Independent Column Approximation (McICA) [*Pincus et al.*, 2003]. The improvement treatment of cloud overlap and sub-grid scale heterogeneity in  $\tau$ , with a retuning of the model, is apparently responsible for the reduction in optically thick cloud. In this example, an improved treatment of the radiative impact of clouds permitted better clouds properties to be simulated in a model that must match the observed radiation budget. A sensitivity study using McICA in the GFDL model [see Figure 4 of *Zhang et al.*, 2005] also shows a noticeable reduction in the amount of optically thick cloud.

In the Hadley Centre models, McICA is not used so other explanations must be sought. The largest reduction in optically thick cloud happened between HadSM3 and HadSM4, with a smaller but still sizeable reduction between HadSM4 and HadGSM1. Between HadSM3 and HadSM4, boundary layer vertical resolution was increased, the *Lock et al.* [2000] boundary layer turbulence parameterization was introduced, as was a sub-grid (in the vertical) treatment of cloud fraction. The possibility for clouds to occur in thinner layers admits the possibility of lower optical depths in stratiform clouds to be simulated (at fixed water content) (M. Webb, personal communication). The vertical resolution of climate models is known to be too coarse to simulate the many stratiform clouds that have geometrical cloud thicknesses smaller than that typical of model layers. Additionally, HadSM4 introduced an improved treatment of mixed-phase cloud microphysics [*Wilson and Ballard*, 1999] which also may be a factor in the reductions of optically thick cloud, particularly at middle-latitudes where a treatment of the Bergeron

process may reduce the amount of super-cooled liquid in deep frontal clouds.

In the Community Atmosphere models, the vertical resolution in the boundary layer was increased and every physical parameterization, except that of deep convection, was changed between CAM4 and CAM5. Thus, all of the explanations above may be playing a role in their reduction of optically thick cloud [Neale *et al.*, 2011]. In particular, the introduction of improved cloud microphysics led to a substantial reduction in liquid water path over middle-latitudes that probably contributes to the reduction of optically thick clouds [Gettelman *et al.*, 2008].

Our evaluation is necessarily incomplete. For example, it is of interest to evaluate other cloud properties, such as liquid and ice water paths, or modes of variability, or how clouds co-vary with environmental parameters including 500 hPa vertical velocity and lower tropospheric stability. Because our analysis requires the use of an ISCCP simulator, our study is limited in the number of models that we can examine, although most major climate models have been included in this study. Evaluation of the limited and less consistently determined cloud information collected from a wider set of climate models is also of interest [Jiang *et al.*, 2012].

One may wonder if there is any connection between improved cloud simulations in climate models and the response to greenhouse gases in the climate changes these model simulate. Previous investigations have found no significant relationship between climate sensitivity and the fidelity of a model simulation in simulating present-day

climate of clouds and precipitation [*Pincus et al.*, 2008]. We note that range of climate sensitivity in CMIP5 models is just as wide as it was in CMIP3 [*Andrews et al.*, 2012], again with the diversity in cloud feedbacks being a leading cause of inter-model spread. This suggests that there is no connection between the global mean cloud feedback and the fidelity with which a model simulates the clouds of the present-day climate. One implication of the reduction of cloud optical depths is that the magnitude of cloud feedbacks resulting from optical depth changes can be substantially larger if the current climate's cloud albedo is not saturated [*Stephens* 2010].

**Acknowledgments.** We acknowledge the World Climate Research Program's Working Group on Coupled Modeling, which is responsible for CMIP, and we thank the climate modeling groups (listed in Tables 1 and 2 of this paper) for producing and making available their model output. For CMIP the U.S. Department of Energy's Program for Climate Model Diagnosis and Intercomparison provides coordinating support and led development of software infrastructure in partnership with the Global Organization for Earth System Science Portals. The efforts of all authors from Lawrence Livermore National Laboratory were supported by the Regional and Global Climate and Earth System Modeling programs of the United States Department of Energy's Office of Science and were performed under the auspices of the United States Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344. RP appreciates support from NASA under grant NNX11AF09G and from NSF under grant AGS 1138394. We thank Ben Sanderson for providing ISCCP simulator

output from the CAM4 slab-ocean model, and Alejandro Bodas-Salcedo for providing additional ISCCP simulator output from the Hadley Center models. We thank Jason Cole, Mark Webb, and Shaocheng Xie for conversations.

## References

Ackerman, T. P., and G. M. Stokes (2003), The Atmospheric Radiation Measurement Program, *Phys. Today*, *56*, 38–44.

Andrews, T., J. M. Gregory, M. J. Webb, and K. E. Taylor (2012), Forcing, feedbacks and climate sensitivity in CMIP5 coupled atmosphere-ocean climate models, *Geophys. Res. Lett.*, *39*, L09712, doi:10.1029/2012GL051607.

Bodas-Salcedo, A., et al. (2011), COSP: Satellite simulation software for model assessment, *Bull. Amer. Meteor. Soc.*, *92*, 1023–1043.

Bony, S., and J.-L. duFresne (2005), Marine boundary layer clouds at the heart of tropical cloud feedback uncertainties in climate models, *Geophys. Res. Lett.*, *32*, L20806, doi:10.1029/2005GL023851.

Bony, S., M. Webb, C. Bretherton, S. Klein, P. Siebesma, G. Tselioudis, and M. Zhang (2011), CFMIP: Towards a better evaluation and understanding of clouds and cloud feedbacks in CMIP5 models, *CLIVAR Exchanges*, *56*, International CLIVAR Project

606 Office, Southampton, United Kingdom, 20-24.

607 Bretherton, C. S. and S. Park (2009), A new moist turbulence parameterization in the  
 608 Community Atmosphere Model, *J. Clim.*, 22, 3422-3448.

609 Collins, W. D. et al. (2006), The formulation and atmospheric simulation of the  
 610 Community Atmosphere Model Version 3 (CAM3), *J. Clim.*, 19, 2144-2161.

611 Collins, W. J. et al. (2008), *Evaluation of the HadGEM2 model*, Met Office Hadley  
 612 Centre Technical Note no. HCTN 74, Met Office, FitzRoy Road, Exeter EX1 3PB,  
 613 United Kingdom.

614 Curry, J. A., W. B. Rossow, D. Randall, and J. L. Schramm (1996), Overview of Arctic  
 615 cloud and radiation characteristics, *J. Clim.*, 9, 1731–1764.

616 Dai, A. (2006), Precipitation characteristics in eighteen coupled climate models, *J. Clim.*,  
 617 19, 4605-4630.

618 Garay, M. J., S. P. de Szoeke, and C. M. Moroney (2008), Comparison of marine  
 619 stratocumulus cloud top heights in the southeastern Pacific retrieved from satellites with  
 620 coincident ship-based observations, *J. Geophys. Res.*, 113, D18204, doi:  
 621 10.1029/2008JD009975.



622 Gates, W. L., et al. (1999), An overview of the results of the Atmospheric Model  
623 Intercomparison Project (AMIP I), *Bull. Amer. Meteor. Soc.*, *80*, 29–55.

624 Gent, P. R. et al. (2011), The Community Climate System Model Version 4, *J. Clim.*, *24*,  
625 4973-4991.

626 Gettelman, A., H. Morrison, and S. J. Ghan (2008), A new two-moment bulk stratiform  
627 cloud microphysics scheme in the Community Atmosphere Model (CAM3), Part II:  
628 Single-column and global results, *J. Clim.*, *21*, 3660-3679.

629 GEWEX Cloud System Science Team (1993), The GEWEX Cloud System Study  
630 (GCSS), *Bull. Amer. Meteor. Soc.*, *74*, 387–399.

631 GFDL GAMDT (2004), The new GFDL global atmosphere and land model AM2/LM2:  
632 Evaluation with prescribed SST simulations, *J. Clim.*, *17*, 4641-4673.

633 Gleckler, P. J., K. E. Taylor, and C. Doutriaux (2008), Performance metrics for climate  
634 models, *J. Geophys. Res.*, *113*, D06104, doi:10.1029/2007JD008972.

635 Hourdin, F. et al. (2006), The LMDZ4 general circulation model: climate performance  
636 and sensitivity to parametrized physics with emphasis on tropical convection, *Clim. Dyn.*,  
637 *27*, 787-813.

638 IPCC (2007), *Climate Change 2007: The Physical Science Basis*, Contribution of  
 639 Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on  
 640 Climate Change, [Solomon, S., D. Qin, M. Manning, Z. Chen, M. Marquis, K. B. Averyt,  
 641 M. Tignor and H. L. Miller (eds.)], Cambridge University Press, Cambridge, United  
 642 Kingdom and New York, NY, USA, 996 pp.

643 Jiang, J. et al. (2012), Evaluation of cloud and water vapor simulations in CMIP5 climate  
 644 models using NASA “A-Train” satellite observations, *J. Geophys. Res.*,  
 645 doi:10.1029/2011JD017237.

646 Kay, J., et al. (2012), Exposing global cloud biases in the Community Atmosphere Model  
 647 (CAM) using satellite observations and their corresponding instrument simulators, *J.*  
 648 *Clim.*, in press.

649 Klein, S. A., and C. Jakob (1999), Validation and sensitivities of frontal clouds simulated  
 650 by the ECMWF model, *Mon. Weather Rev.*, *127*, 2514–2531.

651 Lock, A. P., A. R. Brown, M. R. Bush, G. M. Martin, and R. N. B. Smith (2000), A new  
 652 boundary layer mixing scheme. Part I: Scheme description and single-column model  
 653 tests, *Mon. Wea. Rev.*, *128*, 3187–3199.

654 Ma, C.-C., C. R. Mechoso, A. W. Robertson, and A. Arakawa (1996), Peruvian stratus  
 655 clouds and the tropical pacific circulation: A coupled ocean-atmosphere GCM study, *J.*

656 *Clim.*, 9, 1635–1645.

657 Mace, G. G., Q. Zhang, M. Vaughan, R. Marchand, G. Stephens, C. Trepte, and D.  
658 Winker (2009), A description of hydrometeor layer occurrence statistics derived from the  
659 first year of merged Cloudsat and CALIPSO data, *J. Geophys. Res.*, 114, D00A26,  
660 doi:10.1029/2007JD009755.

661 Mace, G. G., S. Houser, S. Benson, S. A. Klein and Q. Min (2011), Critical evaluation of  
662 the ISCCP simulator using ground-based remote sensing data, *J. Clim.*, 24, 1598–1612.

663 Marchand, R., T. Ackerman, M. Smyth, and W. B. Rossow (2010), A review of cloud top  
664 height and optical depth histograms from MISR, ISCCP, and MODIS, *J. Geophys. Res.*,  
665 115, D16206, doi:10.1029/2009JD013422.

666 Martin, G. M., et al. (2006), The physical properties of the atmosphere in the new Hadley  
667 Centre Global Environmental Model (HadGEM1). Part I: Model description and global  
668 climatology, *J. Clim.*, 19, 1274-1301.

669 McAvaney, B. J., and H. Le Treut (2003), The cloud feedback intercomparison project:  
670 (CFMIP). *CLIVAR Exchanges*, 26, International CLIVAR Project Office, Southampton,  
671 United Kingdom, 1-4.

672 Meehl, G., C. Covey, T. L. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, and R. J.

673 Stouffer and K. E. Taylor (2007), The WCRP CMIP3 multimodel dataset: A new era in  
674 climate change research, *Bull. Amer. Meteor. Soc.*, 88, 13830-1394.

675 Neale, R. B. et al. (2011a), *Description of the NCAR Community Atmosphere Model*  
676 *(CAM5)*, Technical Report NCAR/TN-486+STR, National Center for Atmospheric  
677 Research, Boulder, Colorado, U. S. A., 268 pp.

678 Ogura, T. et al. (2008), Towards understanding cloud response in atmospheric GCMs:  
679 The use of tendency diagnostics, *J. Met. Soc. Japan*, 86, 69-79.

680 Park, S. and C. S. Bretherton (2009), The University of Washington shallow convection  
681 and moist turbulence schemes and their impact on climate simulations with the  
682 Community Atmosphere Model, *J. Clim.*, 22, 3449-3469.

683 Pincus, R., H. W. Barker, and J. Morcrette (2003), A fast, flexible, approximate  
684 technique for computing radiative transfer in inhomogeneous clouds, *J. Geophys. Res.*,  
685 108(D13), 4376, doi:10.1029/2002JD003322.

686 Pincus, R., C. P. Batstone, R. J. P. Hofmann, K. E. Taylor, and P. J. Gleckler (2008),  
687 Evaluating the present-day simulation of clouds, precipitation, and radiation in climate  
688 models, *J. Geophys. Res.*, 113, D14209, doi:10.1029/2007JD009334.

689 Pincus, R., S. Platnick, S. A. Ackerman, R. S. Hemler, R. J. P. Hoffmann (2012),

690 Reconciling simulated and observed views of clouds: MODIS, ISCCP, and the limits of  
 691 instrument simulators, *J. Clim.*, *25*, 4699-4720.

692 Pope, V. D., M. L. Gallani, P. R. Rowntree, and R. A. Stratton (2000), The impact of new  
 693 physical parametrizations in the Hadley Centre climate model – HadAM3, *Clim. Dyn.*,  
 694 *16*, 123-146.

695 Raddatz, T. J. et al. (2007), Will the tropical land biosphere dominate the climate-carbon  
 696 cycle feedback during the twenty first century? *Clim. Dyn.*, *29*, 565-574,  
 697 doi:10.1007/s00382-007-0247-8.

698 Rossow, W. B. and R. A. Schiffer (1991), International Satellite Cloud Climatology  
 699 Project (ISCCP) cloud data products, *Bull. Amer. Meteor. Soc.*, *72*, 2–20.

700 Rossow, W. B. and R. A. Schiffer (1999), Advances in understanding clouds from  
 701 ISCCP, *Bull. Amer. Meteor. Soc.*, *80*, 2261–2288.

702 Slingo, A., and J.-M. Slingo (1988), The response of a general circulation model to cloud  
 703 longwave radiative forcing. I. Introduction and initial experiments, *Quart. J. Roy. Met.*  
 704 *Soc.*, *114*, 1027-1062.

705 Stephens, G. L., et al. (2002), The Cloudsat mission and the A-train, *Bull. Amer. Meteor.*

706 *Soc.*, 83, 1771–1790.

707 Stubenrauch, C., S. Kinne, and the GEWEX Cloud Assessment Team (2009),  
 708 Assessment of global cloud climatologies, *GEWEX Newsletter*, 19, International  
 709 GEWEX Project Office, Silver Spring, Maryland, Unites States of America, 6-7.

710 Stephens, G. (2010), *Is there a missing low-cloud feedback in current climate models?*  
 711 *GEWEX Newsletter*, 20, International GEWEX Project Office, Silver Spring, Maryland,  
 712 United States of America, 5-7.

713 Taylor, K. E., R. J. Stouffer, and G. A. Meehl (2012), An overview of CMIP5 and the  
 714 experimental design, *Bull. Amer. Meteor. Soc.*, 93, 485-498.

715 Thorsen, T. J., Q. Fu, and J. Comstock (2011), Comparison of the CALIPSO satellite and  
 716 ground-based observations of cirrus clouds at the ARM TWP sites, *J. Geophys. Res.*, 116,  
 717 D21203, doi:10.1029/2011JD015970.

718 Voldoire, et al. (2012), The CNRM-CM5.1 global climate model: description and basic  
 719 evaluation, *Clim. Dyn.*, doi:10.1007/s00382-011-1259-y.

720 Webb, M., C. Senior, S. Bony, and J. J. Morcrette (2001), Combining ERBE and ISCCP  
 721 data to assess clouds in the Hadley Centre, ECMWF and LMD atmospheric climate

722 models, *Clim. Dyn.*, 17, 905–922.

723 Williams, K. D. and M. J. Webb (2009), A quantitative performance assessment of cloud  
724 regimes in climate models. *Clim. Dyn.*, 33, 141-157.

725 Wilson, D. R. and S. P. Ballard (1999), A microphysically based precipitation scheme for  
726 the U. K. Meteorological Office Unified Model, *Q. J. Roy. Met. Soc.*, 125, 1607-1636.

727 Winker, D., et al. (2009), Overview of the CALIPSO mission and CALIOP data  
728 processing algorithms, *J. Atmos. Oceanic Technol.*, 26, 2310-2323.

729 Yang, G.-Y. and J. Slingo, (2001), The diurnal cycle in the tropics, *Mon. Wea. Rev.*, 129,  
730 784-801.

731 Yukimoto, S. et al. (2011a), *Meteorological Research Institute – Earth System Model*  
732 *Version 1 (MRI-ESM1): Model Description*, Technical Report #64, Meteorological  
733 Research Institute, Tsukuba-city, Ibaraki 305-0052, Japan, 96 pp.

734 Zelinka, M. D., S. A. Klein and D. L. Hartmann (2012), Computing and partitioning  
735 cloud feedbacks using cloud property histograms. Part I: Cloud radiative kernels, *J.*  
736 *Clim.*, 25, 3715–3735.

737 Zhang, M. H., et al. (2005), Comparing clouds and their seasonal variations in 10

738 atmospheric general circulation models with satellite measurements, *J. Geophys. Res.*,  
739 110, D15S02, doi:10.1029/2004JD005021.

740



741 **Tables**

742 Table 1. CFMIP 1 slab ocean models used in this study.

Model Name	Modeling Center	Reference	Number of Years in Run	Symbol
AGCM4.0	Canadian Centre for Climate Modeling and Analysis	<a href="http://www.ec.gc.ca/ccmac-cccma/">http://www.ec.gc.ca/ccmac-cccma/</a>	20	c4
CCSM3.0	National Center for Atmospheric Research	<i>Collins et al.</i> [2004]	20	n3
GFDL MLM 2.1	NOAA / Geophysical Fluid Dynamics Laboratory	<i>GFDL GAMDT</i> [2004]	20	g
HadGSM1	Met Office Hadley Centre	<i>Martin et al.</i> [2006]	20	h1
HadSM3	Met Office Hadley Centre	<i>Pope et al.</i> [2000]	20	h3
HadSM4	Met Office Hadley Centre	<i>Webb et al.</i> [2001]	20	h4
IPSL CM4	Institut Pierre Simon Laplace	<i>Hourdin et al.</i> [2006]	20	i
MIROC (hisens)	Center for Climate System Research (The University of Tokyo), National Institute for Environmental Studies, and Frontier Research Center for Global Change	<i>Ogura et al.</i> [2008]	5	m3
MIROC (losens)	Center for Climate System Research (The University of Tokyo), National Institute for Environmental Studies, and Frontier Research Center for Global Change	<i>Ogura et al.</i> [2008]	5	m4

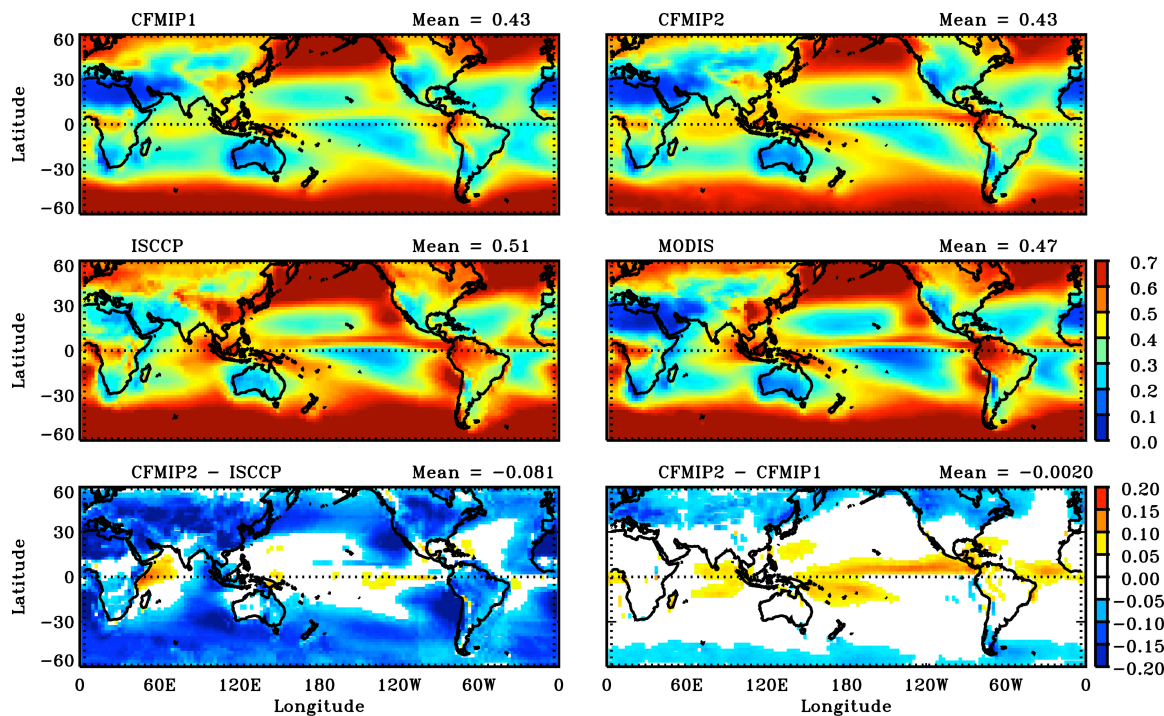
743

744 Table 2. CFMIP 2 AMIP models used in this study.

Model Name	Modeling Center	Reference	Number of Years in Run	Symbol
CAM4	Community Earth System Model Contributors (NSF-DOE-NCAR)	<i>Gent et al. [2004]</i>	10	N4
CAM5	Community Earth System Model Contributors (NSF-DOE-NCAR)	<i>Neale et al. [2011]</i>	10	N5
CanAM4	Canadian Centre for Climate Modeling and Analysis	<a href="http://www.ec.gc.ca/ccmac-cccma/">http://www.ec.gc.ca/ccmac-cccma/</a>	60	C4
CNRM-CM5	Centre National de Recherches Meteorologiques / Centre Europeen de Recherche et Formation Avancees en Calcul Scientifique	<i>Voldoire et al. [2012]</i>	30	Q
HadGEM2A	Hadley Centre for Climate Prediction and Research/Met Office	<i>Collins et al. [2008]</i>	30	H2
MIROC5	Atmosphere and Ocean Research Institute (The University of Tokyo), National Institute for Environmental Studies, and Japan Agency for Marine-Earth Science and Technology	<i>Watanabe et al. [2010]</i>	30	M5
MPI-ESM-LR	Max Planck Institute for Meteorology	<i>Raddatz et al. [2007]</i>	30	P
MRI-CGCM3	Meteorological Research Institute	<i>Yukimoto et al. [2011]</i>	30	R

745

746 **Figures**



748 Figure 1. Total cloud amount ( $\tau > 1.3$ ) from CFMIP1 and CFMIP2 multi-model means,  
749 ISCCP and MODIS observations, and the difference of CFMIP2 multi-model mean to the  
750 ISCCP and CFMIP1 multi-model mean.

751

752

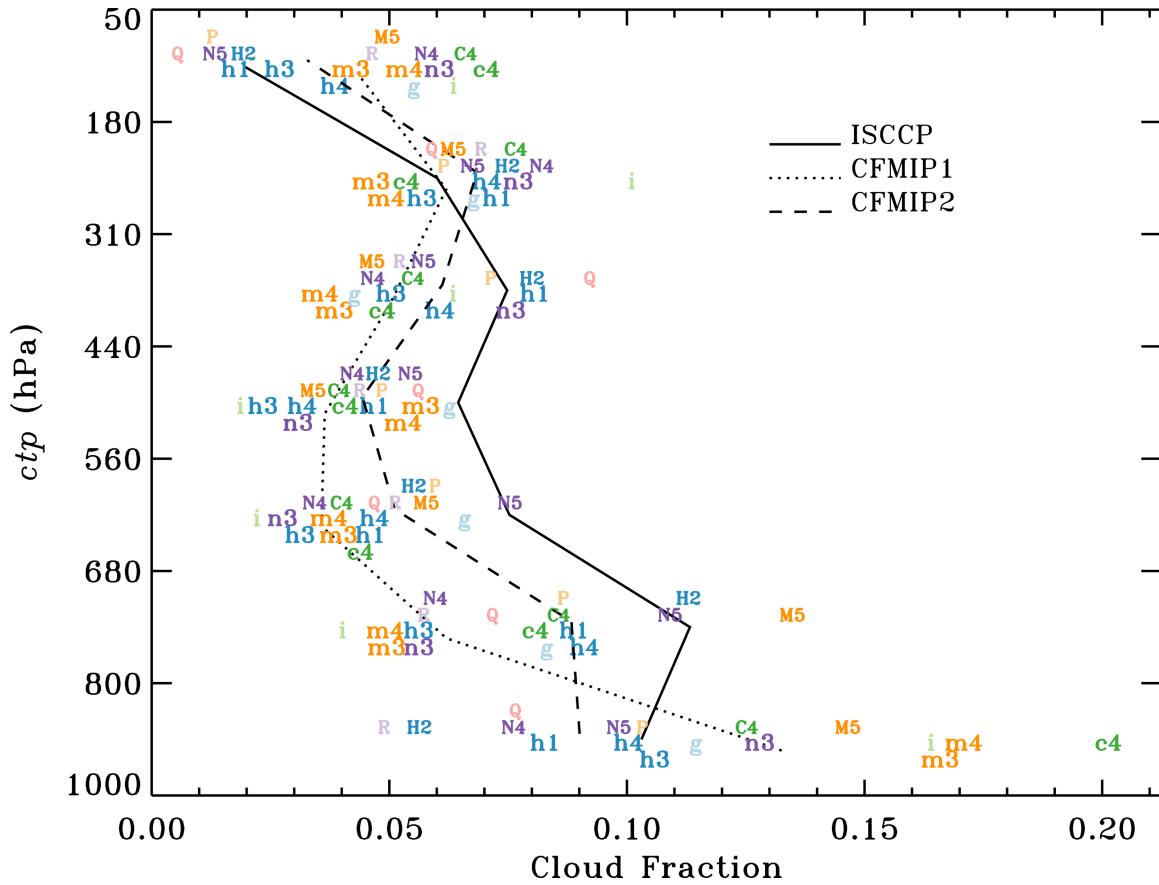


Figure 2. Fractional area in the domain 60°S - 60°N covered by clouds as a function of cloud-top pressure from models and ISCCP observations. CFMIP1 (2) ensemble means are plotted with a dotted (dashed) line. The area is computed only for clouds with  $\tau > 1.3$ . The symbol key for models is provided in Tables 1 and 2.



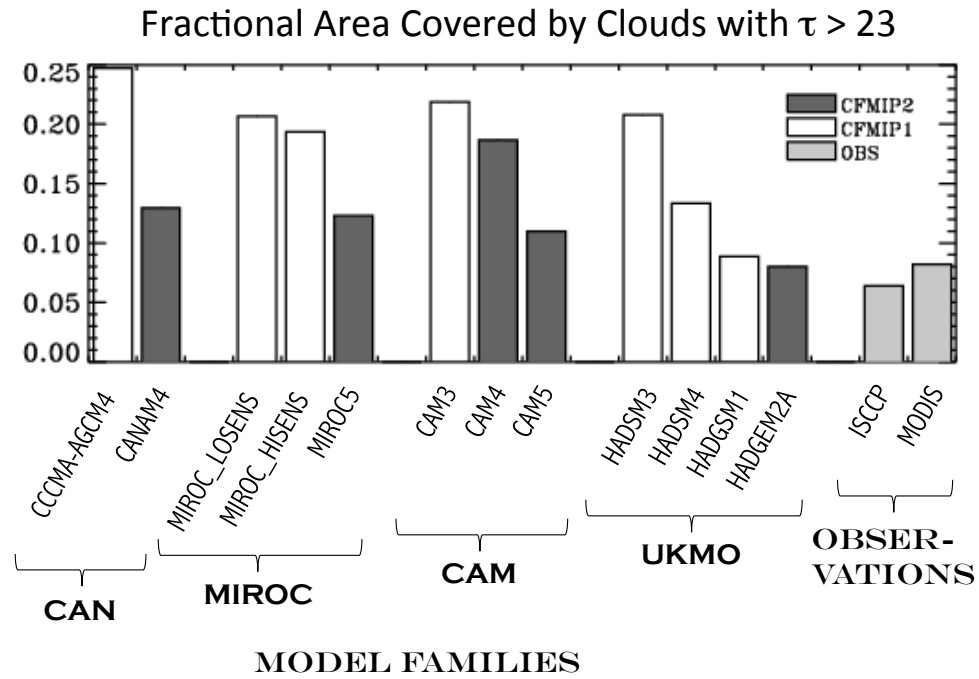


Figure 4. Fractional area in the domain 60°S - 60°N covered by clouds with  $\tau > 23$  for selected model families and observations. Models are plotted so as to illustrate progress in reducing the overestimate of optically thick cloud over time by ordering models from earliest to latest (left to right) within families.

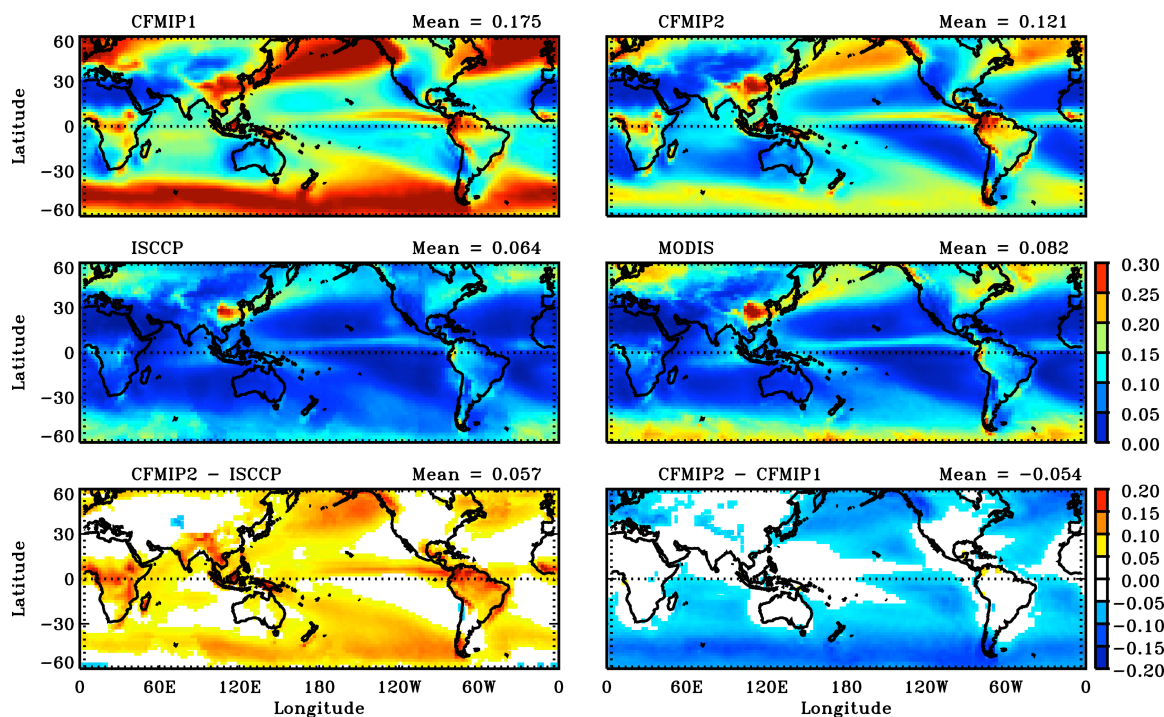
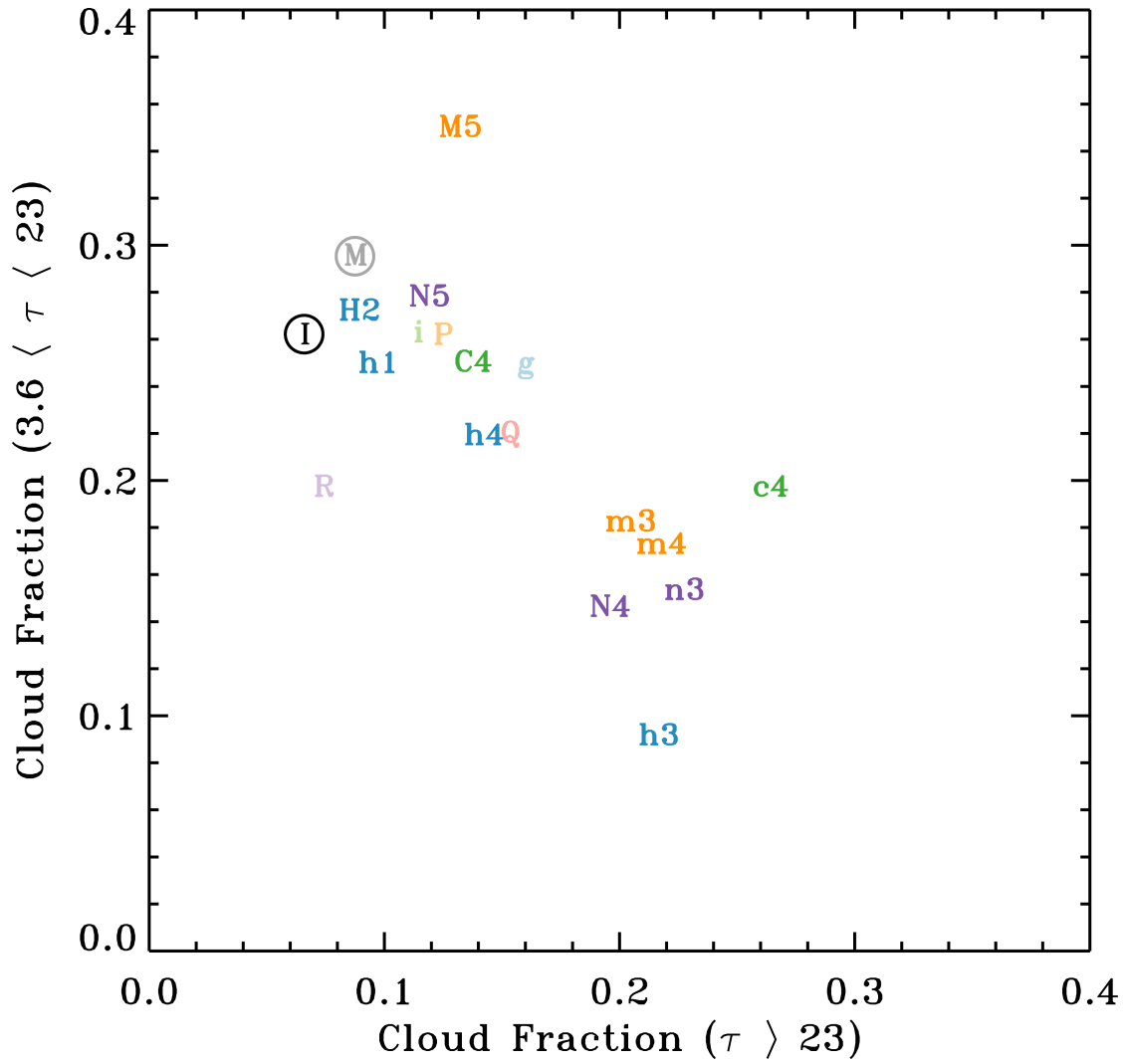


Figure 5. Fractional area covered by optically thick clouds ( $\tau > 23$ ) from CFMIP1 and CFMIP2 multi-model means, ISCCP and MODIS observations, and the difference of the CFMIP2 multi-model mean to ISCCP and the CFMIP1 multi-model mean.

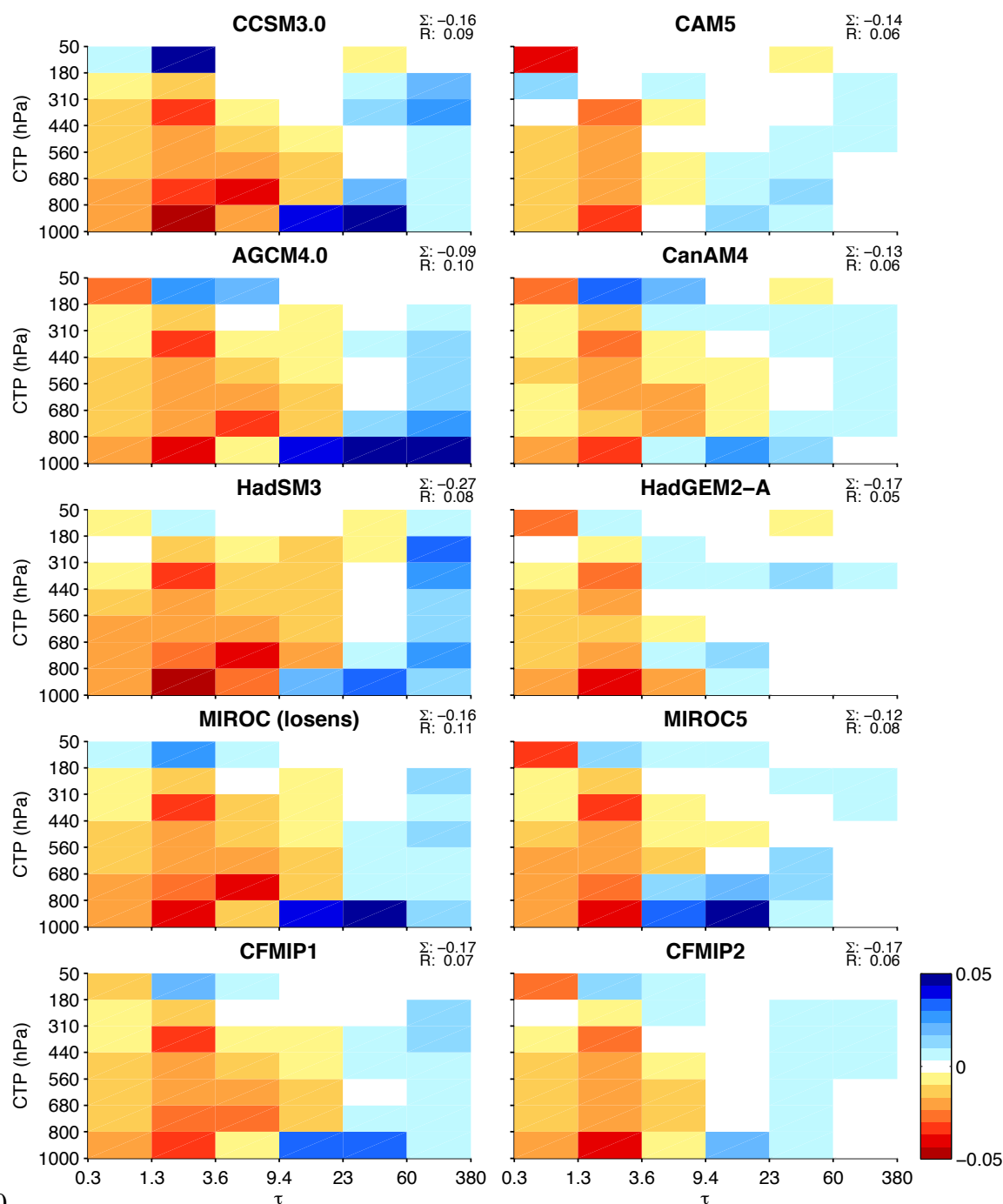


774

775 Figure 6. Scatterplot of the fractional area in the domain 60°S - 60°N covered by clouds  
 776 with  $\tau > 23$  and clouds with  $3.6 < \tau < 23$ . Observations from MODIS and ISCCP are  
 777 represented by “M” and “I”, respectively. The symbol key for models is provided in  
 778 Tables 1 and 2.

779





780

781 Figure 7. Area-averaged biases in the domain 60°S - 60°N with respect to ISCCP observations of fractional area  
 782 covered by clouds in bins of cloud-top pressure and optical depth. Results are plotted for the 4 model families in  
 783 which we track progress and the ensemble mean. Models are ordered with the oldest models on the left and the  
 784 newest models on the right. The sum of the histogram and the range (maximum minus minimum value in the  
 785 histogram) are shown in the title of each panel. Positive values indicate model overestimates relative to  
 786 observations. The fact that the recent models have fewer bins with color as well as reduced intensity in the bins  
 787 with color indicates improvements with time.

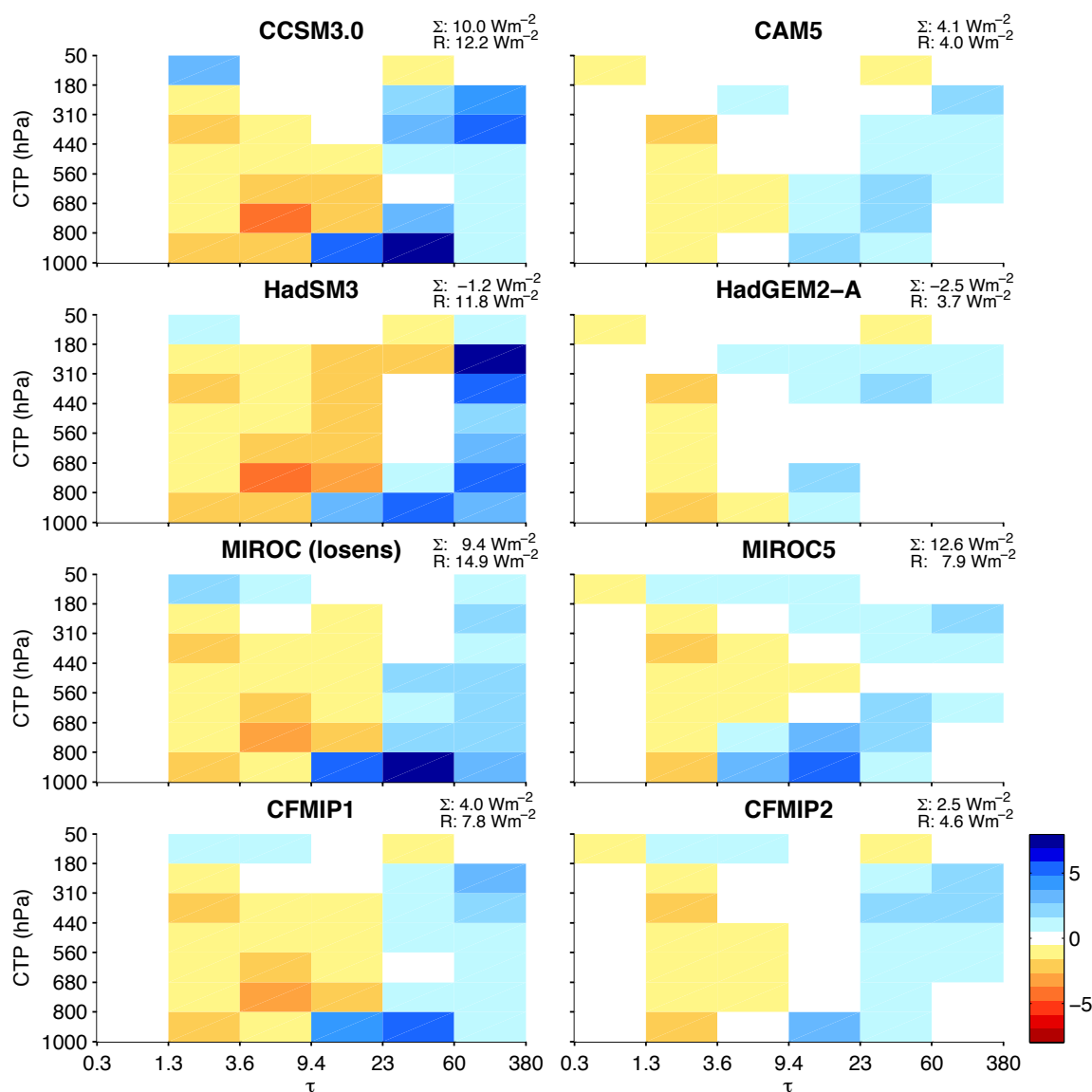


Figure 8. As in Figure 7, but for the contributions to shortwave radiation reflected to space by clouds stratified into bins of cloud-top pressure and optical depth. Positive values indicate a bias towards too much reflected radiation due to a positive bias in cloud amount.

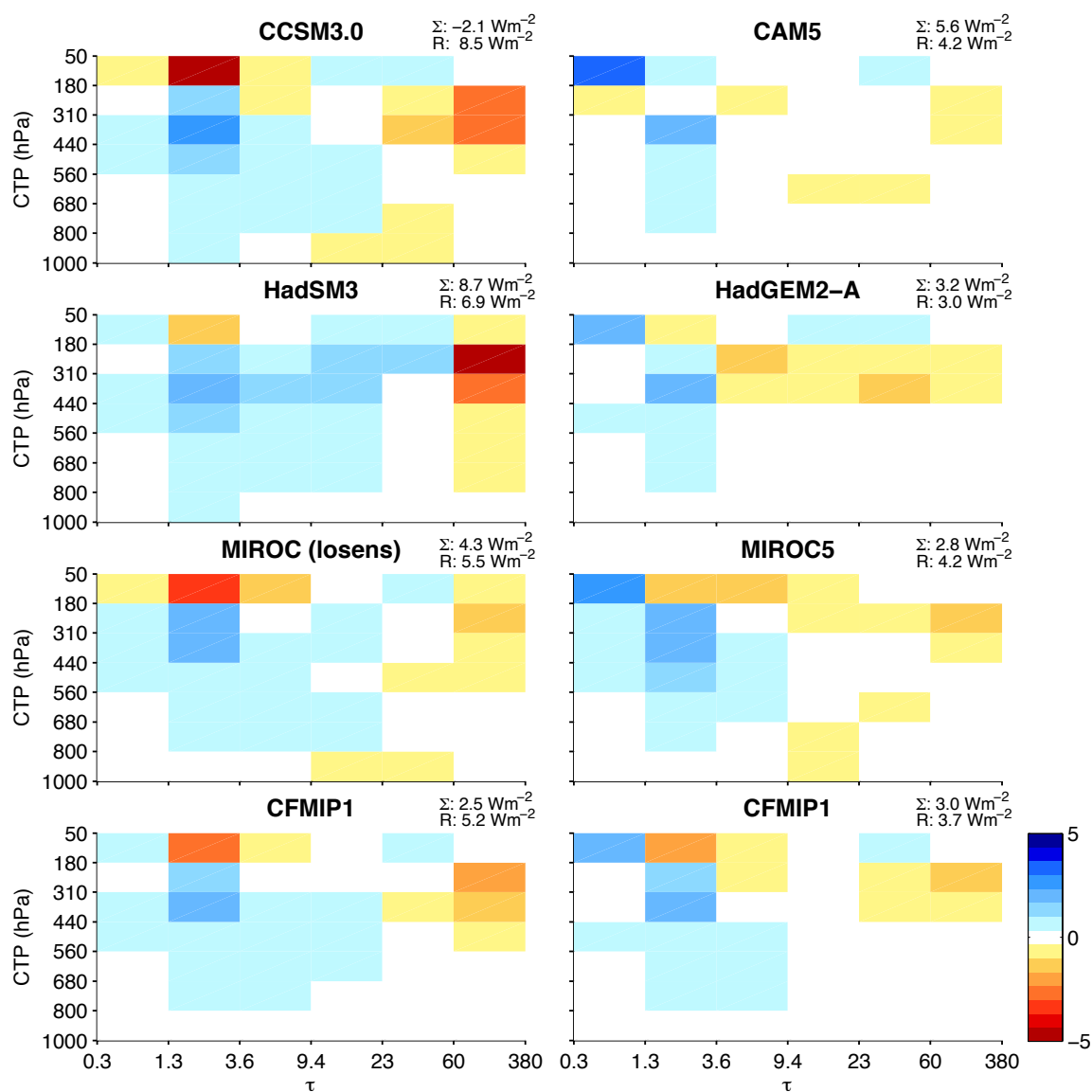
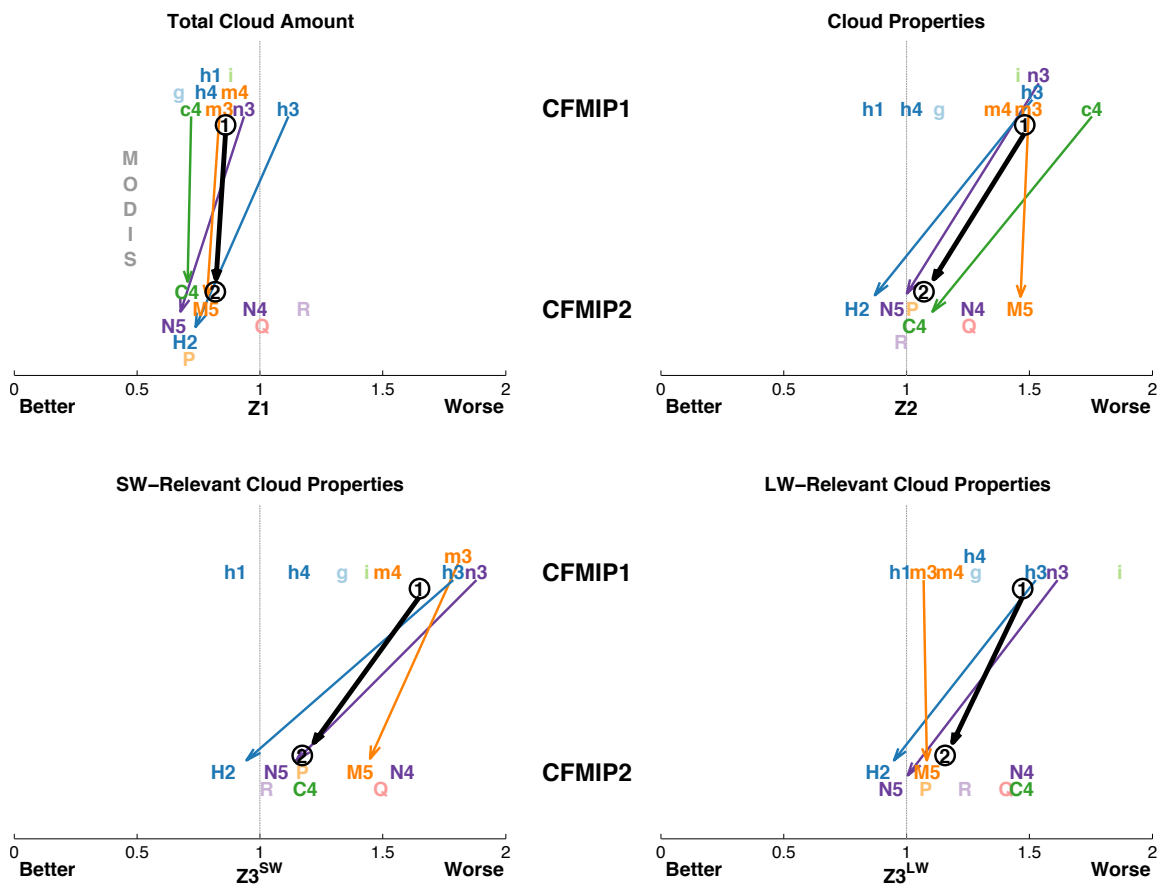


Figure 9. As in Figure 7, but for the contributions to reductions of outgoing longwave radiation (relative to clear-sky) by clouds stratified into bins of cloud-top pressure and optical depth. Positive values indicate a bias towards too much longwave radiation emitted to space due to a negative bias in cloud amount.



800

801 Figure 10. Scalar measures of fidelity of CFMIP model simulations in reproducing the  
 802 space-time distribution of several cloud measures, with greater fidelity indicated by lower  
 803 Z values. Z<sub>1</sub> measures fidelity in simulating total cloud amount, whereas Z<sub>2</sub> measures  
 804 fidelity in simulating cloud-top pressure and optical depth in different categories of  
 805 optically intermediate and thick clouds at high, middle, and low-levels of the atmosphere.  
 806 Z<sub>3</sub> measures the impacts on top-of-atmosphere shortwave (lower left) and longwave  
 807 (lower right) radiation in the same categories measured by Z<sub>2</sub>. Models are stratified  
 808 vertically into the two ensembles and are plotted according to the symbol key in Tables 1  
 809 and 2. For the modeling centers in which we can track progress, the arrow connects the  
 810 oldest model in the family (arrow base) to the most recent model (arrow tip). The thick  
 811 black arrow connects the average measure of CFMIP1 models (arrow base) to that of  
 812 CFMIP2 models (arrow tip). Arrows pointing to the left indicate improvements with  
 813 time.